

Balanced Vertices in Trees and a Simpler Algorithm to Compute the Genomic Distance[☆]

Péter L. Erdős^{a,1}, Lajos Soukup^{a,2}, Jens Stoye^{b,*}

^aAlfréd Rényi Institute of Mathematics, Budapest, P.O. Box 127, H-1364 Hungary

^bUniversität Bielefeld, Technische Fakultät, AG Genominformatik, 33594 Bielefeld, Germany

Abstract

This paper provides a short and transparent solution for the covering cost of white–grey trees which play a crucial role in the algorithm of Bergeron *et al.* to compute the rearrangement distance between two multichromosomal genomes in linear time (*Theor. Comput. Sci.*, 410:5300–5316, 2009). In the process it introduces a new *center* notion for trees, which seems to be interesting on its own.

Keywords: comparative genomics, genome rearrangement, combinatorics on trees

1. Introduction

Computational comparative genomics is a subdiscipline of computational biology in which the relationships between two or more genomes are studied by computational means. A highly relevant question in this field is the calculation of the minimum number of rearrangement operations (reversals, translocations, fusions and fissions) that are necessary to transform one given genome into another one, the so-called *genome rearrangement problem* [1].

The *white–grey tree cover problem* studied in this paper (formally defined in Section 2) arises as a subproblem of the genome rearrangement problem, and so far only an unsatisfactory (and not self-contained) solution exists [2]. The main goal of this paper is to give a short solution of the problem and to correct some omissions and discrepancies of the original formulation. (In Section 4 we point out some cases where the original formulation fails.) Moreover, it gives rise to a combinatorial problem on trees, detailed in Section 3, that seems to be interesting on its own. Since one of our main concerns here is brevity, we (usually) don’t give detailed proofs of easy facts, which are not essential for our main goal.

[☆]This research was carried out when the first and second authors visited the third author at Bielefeld University with the support of the Hungarian Bioinformatics MTKD-CT-2006-042794, Marie Curie Host Fellowships for Transfer of Knowledge.

*Corresponding author

Email addresses: elp@renyi.hu (Péter L. Erdős), soukup@renyi.hu (Lajos Soukup), stoye@techfak.uni-bielefeld.de (Jens Stoye)

¹Partly supported by Hungarian NSF, under contract Nos. AT048826 and K68262.

²Partly supported by Hungarian NSF, under contract Nos. K61600 and K68262.

2. Problem definition

A *white–grey tree* is a rooted tree with (white or grey) colored and uncolored vertices. The **root** is uncolored, some children of the **root** are grey (some of them can be leaves), all leaves which are not children of the **root** are white. All uncolored vertices (with the possible exception of the root) are branching points.

A system of paths in a white–grey tree is a *colored cover* if:

- (i) Each path has colored endpoints. One vertex alone may constitute a path.
- (ii) Each colored vertex is covered with path(s).

The *cost* of a path P is denoted by $\text{cost}(P)$ and is defined as follows:

- (i) P is *short* if it has exactly one vertex. Then $\text{cost}(P) = 1$.
- (ii) P is *grey* if its endpoints are grey vertices (then the third vertex is the uncolored **root**). Then $\text{cost}(P) = 1$.
- (iii) P is *long* otherwise. Its cost is $\text{cost}(P) = 2$.

Definition 1. The *cost of a path system \mathcal{P}* is the sum of the individual costs: $\text{cost}(\mathcal{P}) := \sum_{P \in \mathcal{P}} \text{cost}(P)$. A colored cover \mathcal{P} is an *optimal* one for a given white–grey tree T if it has minimal cost among all possible colored covers, denoted by $\text{cost}(T)$.

Problem 2 (White–grey tree cover problem). Given a white–grey tree T , compute $\text{cost}(T)$.

The main result of this paper is a simple way to calculate the exact cost of an optimal cover. We are not quite ready to formalize the main result (without some further observations it would require a detailed case analysis), but we mention here a well known fact [1]:

Lemma 3. *Let T be a white–grey tree with w white and g grey leaves, then:*

$$w + \lceil g/2 \rceil \leq \text{cost}(T) \leq w + \lceil g/2 \rceil + 1.$$

□

3. Balanced vertices in trees

In this section we prove a useful tool for (unrooted) trees which seems to be interesting in its own. In tree T' denote by $P_{u,v}$ the unique path between vertices u and v .

Theorem 4. *Let T' be a tree with $2n$ leaves. Then there exists a vertex $v \in V(T')$ and a bijection among the leaves $\alpha : \mathcal{L} \rightarrow \mathcal{L}$ such that the path system $P_{\ell, \alpha(\ell)}$ (where $\ell \in \mathcal{L}$) covers each vertex in T' and all paths contain v .*

We offer here two proofs. One gives a very simple algorithm to construct such a cover, but it clearly cannot provide all possible solutions. The second proof is based on a necessary and sufficient reformulation of the statement.

FIRST PROOF. Consider an embedding of our tree into the plane and enumerate the leaves in a counter clockwise fashion. One way to obtain such a numbering is to fix a leaf as a root and take the left-to-right, depth first traversal of the tree which conforms with the embedding.

Now we define our bijection with the formula $\alpha : \ell_i \mapsto \ell_{i+n} \bmod 2n$. Considering any two such paths, their endpoints alternate along the circle which contains the leaves in increasing order. Therefore these two paths clearly intersect each other.

As it is well known (its very first proof is due to Gyárfás and Lehel, [3]) if (in a tree) a set of paths does not contain two disjoint paths, then all the paths share a common vertex v . And because these paths connect v to the leaves, they cover all edges of the tree. \square

If T' is a fully balanced tree, then no matter what is the embedding in the previous proof, two close leaves will be paired with two close leaves. Therefore there are clearly solutions which cannot be obtained with the previous method. In the remaining part of this section we sketch a proof which is able to find all possible solutions:

SECOND PROOF. For each vertex-edge pair (v, e) denote by $\delta(v, e)$ the number of leaves ℓ in T' such that $P_{v, \ell}$ contains the edge e (where $v \in V(T')$, $e \in E(T')$ and $v \in e$). Furthermore, denote by $E(v)$ the set of edges that contain v .

In the configuration required by Theorem 4, vertex v clearly satisfies the inequality:

$$\forall e \in E(v) \quad : \quad \delta(v, e) \leq \sum \{ \delta(v, f) : f \in E(v), f \neq e \}. \quad (1)$$

Such a vertex $v \in T'$ is called a *balanced* vertex. (If a vertex-edge pair does not satisfy this inequality, then the pair is called *oversaturated*.) As a matter of fact, this property is just equivalent to the existence of the required cover:

Lemma 5. *Let T' be a tree with $2n$ leaves, $n \geq 1$. Then for any balanced vertex v there exists a bijection α such that the paths $P_{\ell, \alpha(\ell)}$ cover all edges, and all paths contain v .*

The easy proof is left to the diligent reader. (One can argue, for example, with mathematical induction.) A balanced vertex in a tree is similar to the well-known notion *center* of the tree, but while a center is usually (almost) unique, there may be several balanced vertices.

The following observation completes our second proof of Theorem 4.

Lemma 6. *Any tree T' with an even number of leaves contains a balanced vertex.*

PROOF. (Sketch) Assume that a particular vertex v is not balanced. Then there exists an edge $e \in E(v)$ such that the pair (v, e) is oversaturated. We repeat the process with the other end of that edge. If this vertex is not balanced again, then it will provide another oversaturated pair. The finiteness of the graph finishes the proof of the Lemma and this also completes the second proof of Theorem 4. \square

The flexibility in the pairing algorithm clearly can provide any possible bijection α . It is also interesting to recognize that one can find a suitable balanced vertex quickly:

Lemma 7. *Let T' be a tree with $2n$ leaves. Then there is a linear (in the number of leaves) time algorithm to find a balanced vertex in T' .*

This proof is left to the reader again. A simple dynamic programming algorithm suffices. \square

4. Optimal colored covers

We are ready to determine the cost of an optimal cover for the white–grey tree T . We say a path in the cover is a *mixed* path if it contains at least two colored vertices, exactly one of the colored vertices is a grey leaf. We will use the notation T_w for the subtree derived from T by deleting

all grey leaves and their edges, and the **root** if it would become a leaf. Furthermore for a path P in T we will use the notation $P \upharpoonright T_w$ to denote the *trace* of P on T_w , i.e. the restriction of P to the nodes of T_w with the extra condition that in the truncated path we delete the starting (if any) uncolored vertices. We extend this notation to the trace of a path system \mathcal{P} , $\mathcal{P} \upharpoonright T_w$.

Our general strategy to determine an optimal colored cover is to build it from an optimal colored cover of the subtree T_w . To do that we are going to exploit certain properties – described in the following result – of optimal solutions having a minimum number of mixed paths.

Theorem 8. *Every white–grey tree T has an optimal colored cover \mathcal{P} such that*

- (1) \mathcal{P} contains at most 2 mixed paths,
- (2) $\mathcal{P} \upharpoonright T_w$ is an optimal cover of T_w ,
- (3) for each mixed path $P \in \mathcal{P}$, $\text{cost}(P) = \text{cost}(P \upharpoonright T_w)$ and so $P \upharpoonright T_w$ is either a long path, or a short path consisting of a single grey leaf.

PROOF. (1) Let \mathcal{P} be an optimal cover with a minimum number of mixed paths. Assume on the contrary that \mathcal{P} contains three mixed paths: P_1, P_2 and P_3 , where P_i is a path from the grey leaf g_i to the colored vertex $c_i \in T_w$. (If two paths cover the same grey leaf then deleting that leaf from one of the paths decreases the number of the mixed paths in the cover. So we may assume that the grey leaves are pairwise distinct.)

Let the path P be the intersection of the paths P_1, P_2, P_3 . Clearly P is a path from the **root** to some $c \in T_w$. (It is clear that c may be the **root** itself). Since c is the “last” point of the intersection, we can assume that the unique sub-paths $P_{c_1,c}$ and $P_{c_2,c}$ are edge disjoint (and of course vertex-disjoint except vertex c). Then replace the paths $\{P_1, P_2, P_3\}$ in \mathcal{P} with the paths $\{P_{g_1,g_2}, P_{c_1,c_2}, P_3\}$ to obtain a path cover \mathcal{P}' . But $\text{cost}(\mathcal{P}') \leq \text{cost}(\mathcal{P})$ and \mathcal{P}' contains less mixed paths than \mathcal{P} – a contradiction.

(2) So we have an optimal cover which contains at most two mixed paths. If its trace is not optimal then consider the following cover Q : cover T_w optimally (this has cost at least 1 smaller than the trace of the original cover had), keep the paths from \mathcal{P} which do not intersect T_w and finally cover the (at most two) grey vertices that were covered by the mixed path(s) with a path whose cost is 1. Then the cost of Q is less than or equal to the cost of \mathcal{P} , and Q does not contain any mixed path.

(3) Assume that $2 = \text{cost}(P) > \text{cost}(P \upharpoonright T_w) = 1$ for a mixed path $P \in \mathcal{P}$. The restriction $P \upharpoonright T_w$ should be a short white path covering vertex u , and P is a path $P_{u,g} = (u, \text{root}, g)$ for a grey leaf g . Replacing the path P with two short paths covering u and g resp. keeps the cost of the cover, but decreases the number of mixed paths. \square

A cover \mathcal{P} is *nice* iff it satisfies the requirements of Theorem 8. Let P be a path in T_w . We say that path P is *free* iff P can be extended to path P' such that P' contains a grey leaf while $\text{cost}(P) = \text{cost}(P')$ holds. Theorem 8 implies the following statement:

Lemma 9. *Assume that T is a white–grey tree which has g grey leaves.*

$$\text{cost}(T) = \text{cost}(T_w) + \max \left\{ 0, \left\lceil \frac{g-f}{2} \right\rceil \right\},$$

where f is the maximal number of free paths in a nice optimal cover of T_w . \square

Next we should solve the white–grey tree cover problem for the subtree T_w . Therefore we first solve the problem for trees where (essentially) all leaves are white. In what comes, we will say a leaf is *short* if it is adjacent to a branching vertex.

Lemma 10. *Let T' be a white–grey tree with w colored leaves but without a grey vertex or with exactly one grey leaf. Then the minimal cost of a colored cover is:*

$$\text{cost}(T') = \begin{cases} w + 1, & \text{if } w \text{ is odd and there is no short leaf;} \\ w, & \text{otherwise.} \end{cases} \quad (2)$$

PROOF. Since we have at most one grey leaf, we can not use a “cheap” grey-grey path to cover it. So we can change the color of that vertex into white without changing the cost of the tree and thus assume that all leaves are white.

If the number of leaves is even, then the result is a direct consequence of Lemma 3 and Theorem 4. If the number of leaves is odd, but there is a short leaf, then we cover that leaf with a short path. Deleting it from the tree we are back to the previous case.

Finally assume that w is odd but $\text{cost}(T') = w$. Then each leaf is covered once in an optimal cover, and one of them is covered by a short path. If this leaf is not a short one, however, then its colored neighbor is not covered, a contradiction. For simplicity we fix: in this case the constructed optimal cover contains a long path which does not cover any branching vertex. This path will be called a *half-path*. \square

Let’s remark that Lemma 10 for white–only trees is certainly not new: actually it was proved as early as 1995 (see [1]). But the consideration of more general white–grey trees raises several problematic issues. One of them is that in the literature, known to the authors, grey vertices which are not leaves have not been studied. However, the white–grey trees are constructed in connection with the genome rearrangement problem ([2]) and grey vertices can appear in non-leaf positions.

Another problem that paper [2] fails to determine is the exact cost of a minimal colored cover for some cases. Here we give only one of them. (The references relate to the relevant sections of that paper.) Assume that the **root** of T has two neighbors: one is a grey leaf ($g = 1$), and the other one is a branching vertex. Furthermore assume that w is odd, and no white leaf is short. Then we are in the scope of Theorem 5 of [2]. Since g is odd and T_c is a *fortress* or *junior fortress*, we are to apply the case “otherwise” of Theorem 5. That formula now gives: $\text{cost}(T) = w + \lceil g/2 \rceil + 1 = w + 2$ while the proper cost is only $w + \lceil g/2 \rceil = w + 1$.

Before we give our main result we introduce one more notion: when among the children of the **root** there is exactly one child that is not a grey leaf, then the (colored) vertices between the **root** and the first branching point are called *dangerous*.

Theorem 11. *Let T be a white–grey tree with g grey and w white leaves. Let T_w be derived from T by deleting the grey leaves (and the **root** if it would become a leaf).*

(1) *If T does not have any dangerous vertex then*

$$\text{cost}(T) = \begin{cases} w + 1 + \left\lceil \frac{g-1}{2} \right\rceil, & \text{if } w \text{ is odd and there is no short leaf in } T_w; \\ w + \left\lceil \frac{g}{2} \right\rceil, & \text{otherwise.} \end{cases}$$

(2) If T has some dangerous vertices (and T_w has $(w + 1)$ leaves) then

$$\text{cost}(T) = \begin{cases} (w + 1) + 1 + \max \left\{ 0, \left\lceil \frac{g-2}{2} \right\rceil \right\}, & \text{if } (w + 1) \text{ is odd and there is no short leaf in } T_w; \\ (w + 1) + \left\lceil \frac{g}{2} \right\rceil, & \text{if } \begin{cases} (w + 1) \text{ is odd, there is only one short leaf in } T_w, \\ \text{and that leaf is white and dangerous;} \end{cases} \\ (w + 1) + \left\lceil \frac{g-1}{2} \right\rceil, & \text{otherwise.} \end{cases}$$

PROOF. (1) Assume that w is odd and there is no short leaf in T_w . Then, due to Lemma 10, $\text{cost}(T_w) = w + 1$ and the half-path of the solution is clearly free, so $f = 1$. Otherwise $\text{cost}(T_w) = w$, but we have no free path at all, so $f = 0$. In both cases apply Lemma 9.

(2) **Case 1:** Assume that $(w + 1)$ is odd and there is no short leaf in T_w . Then $\text{cost}(T_w) = (w + 1) + 1$. In the derived optimal cover \mathcal{P}_w of T_w a leaf-leaf type long path P covers the dangerous vertices. Then the path P and the half-path of \mathcal{P}_w are free, so $f = 2$. Then apply Lemma 9

Case 2: Assume now that $(w + 1)$ is odd, there is only one short leaf ℓ in T_w , and that leaf is white and dangerous. Then $\text{cost}(T_w) = w + 1$. Moreover, an optimal cover of \mathcal{P}_w should contain the short (non-free) path covering ℓ , and there is no other free path, thus $f = 0$. Lemma 9 finishes the case.

Case 3: The “otherwise” cases: Assume first that $(w + 1)$ is odd, there is only one short leaf ℓ , and that vertex is a grey (therefore also a dangerous) vertex. Then $\text{cost}(T_w) = w + 1$. Moreover, an optimal cover of T_w should contain the short path covering ℓ . But ℓ is grey, so the path covering ℓ is free. Thus $f = 1$. Then apply Lemma 9.

Assume now that $(w + 1)$ is odd, and there is a short leaf ℓ which is not dangerous. Then there is an optimal cover of T_w in which the dangerous vertices are covered by a long path P . Then P is free, so $f \geq 1$. Since $\text{cost}(T_w) = (w + 1)$, in an optimal cover \mathcal{P}_w of T_w there is only one long path which is free because all long paths in \mathcal{P}_w contain two leaves. So $f \leq 1$, i.e. $f = 1$. Now apply Lemma 9.

Assume finally that $(w + 1)$ is even. Then $\text{cost}(T_w) = (w + 1)$, in an optimal cover of T_w there is only one long path which is free because all long paths in an optimal cover should contain two leaves. So $f \leq 1$. However, there is an optimal cover of T_w containing a long path which covers the dangerous vertices. So $f \geq 1$. Thus $f = 1$. Now apply Lemma 9. \square

- [1] S. Hannenhalli, P. A. Pevzner, Transforming men into mice (polynomial algorithm for genomic distance problem), in: Proc. 36th Annu. Symp. Found. Comput. Sci., FOCS 1995, IEEE Press, 1995, pp. 581–592.
- [2] A. Bergeron, J. Mixtacki, J. Stoye, A new linear time algorithm to compute the genomic distance via the double cut and join distance, Theor. Comput. Sci. 410 (2009) 5300–5316.
- [3] A. Gyárfás, J. Lehel, A helly-type problem in trees, in: P. Erdős, A. Rényi, V. T. Sós (Eds.), Combinatorial Theory and its Applications II, Vol. 4 of Colloquia Mathematica Societatis János Bolyai, North Holland, 1970, pp. 571–584, (Balatonfüred, 1969).